

1 Dynamic action recognition based on Dynemes and 2 Extreme Learning Machine

3 Alexandros Iosifidis, Anastasios Tefas and Ioannis Pitas

4 *Department of Informatics, Aristotle University of Thessaloniki*
5 *Thessaloniki 54124, Greece Tel,Fax: +30-2310996304*

6 {aiosif,tefas,pitas}@aiia.csd.auth.gr

7 **Abstract**

In this paper, we propose a novel method that performs dynamic action classification by exploiting the effectiveness of the Extreme Learning Machine (ELM) algorithm for single hidden layer feedforward neural networks training. It involves data grouping and ELM based data projection in multiple levels. Given a test action instance, a neural network is trained by using labeled action instances forming the groups that reside to the test sample's neighborhood. The action instances involved in this procedure are, subsequently, mapped to a new feature space, determined by the trained network outputs. This procedure is performed multiple times, which are determined by the test action instance at hand, until only a single class is retained. Experimental results denote the effectiveness of the dynamic classification approach, compared to the static one, as well as the effectiveness of the ELM in the proposed dynamic classification setting.

8 *Keywords:* Activity recognition; Dynamic classification; Fuzzy Vector
9 Quantization; Extreme Learning Machine

10 **1. Introduction**

11 Human action recognition is a very active research field finding application in
12 many important tasks, such as visual surveillance [1], human-computer interaction
13 [2], augmented reality [3] and semantic video annotation [4]. Actions are usually
14 described by using either features based on optical flow [5], or features devised
15 mainly for action representation [6]. Although the use of such features leads to
16 satisfactory action recognition results, their computation is expensive. Thus, when
17 fast operation is important, action recognition methods should employ simpler ac-
18 tion representations. Neurobiological studies [7] have concluded that the human
19 brain can perceive actions by observing only the human body poses during action
20 execution. Thus, actions can be described as sequences of consecutive human
21 body poses, in terms of human body silhouettes [8, 9]. After describing actions,
22 most methods in the literature exploit supervised machine learning techniques for
23 action class representation and classification of new, unknown, action instances.
24 Such techniques require a training phase, where labeled data are used in order
25 to determine the system parameters. For example, in Artificial Neural Networks
26 (ANNs) based data classification [10], training data are employed in order to de-

27 termine the neurons' weights and in Linear Discriminant Analysis (LDA) based
28 data projection [11], labeled data are used in order to determine a mapping to a
29 lower dimensional feature space for class representation and data classification.
30 Traditionally, the training phase is performed offline by using the entire training
31 set.

32 Action recognition is not an easy task, mainly due to the fact that there is not
33 a formal description of actions. Variations between different action realizations
34 resulted from different action execution styles and different human body sizes be-
35 tween persons result to high intra- and, possibly, low inter-action class variations.
36 This is why person specific classification schemes have been recently investigated
37 for action recognition [12]. The main idea in these classification schemes is to
38 focus the classification problem on each individual person. That is, action recog-
39 nition is performed by a classifier which has been trained by using action instances
40 of the person under consideration. Following this approach, the above mentioned
41 issues are effectively addressed leading to high action classification rates. How-
42 ever, the application of such classification schemes is limited, since, in order to
43 operate properly, a person should have been recorded and trained before recogni-
44 tion. In different cases their performance will probably decrease.

45 An alternative choice could be the use of dynamic action classification schemes.

Dynamic classification, involves a system parameters adaptation procedure based either on the training set structure, or on the test data to be classified. Following this approach, several dynamic classification schemes have been proposed. Wright et. al. [13] proposed a dynamic classification scheme exploiting sparsity constraints. A given test sample is involved in a class independent regression procedure exploiting a codebook containing all the available labeled samples. Multiple reconstruction samples are, subsequently, produced by employing the labeled samples belonging to each class independently. Finally, the test sample is classified to the class providing the minimum reconstruction error. Tang et. al. [14] proposed the Dynamic Committee Machine (DCM), which employs five state-of-the-art classifiers (experts). A test sample is introduced to all the five classifiers and five classification results are produced. The dynamic nature of DCM is based on the adopted fusion strategy, where the experts' weights are modified depending on the corresponding test sample. Kyperountas et. al. proposed a dynamic classification scheme involving an iterative grouping procedure combined with LDA-based data classification [15]. The iterative procedure used in order to determine the optimal training set for LDA based data classification is intuitive and effective. However, the LDA based classification approach in this setting has two disadvantages: a) It sets the assumption of linear class separability. As it will be shown

65 in Section 5, this assumption is not met for action classes. b) The use of a small
66 number of training data, compared to the training data dimensionality, leads to the
67 small sample size problem [16]. In order to address this issue, Kyperountas et.
68 al. employed an LDA variant proposed in [17], in which a regularization param-
69 eter should be a priori known and, thus, an offline training procedure is required.
70 Finally, c) by using training data belonging to only two (or three) classes, LDA
71 projection provides an one- (or two-) dimensional feature space, where classes
72 discrimination may not be captured properly, especially for linear classification
73 models.

74 In order to take into account the non linear nature of action classes, non linear
75 classification methods should be employed. ANNs could be a good choice, since
76 they have proven their effectiveness in a wide range of challenging classification
77 problems. Among them, single hidden layer feedforward networks (SLFNs) have
78 been widely used due to their ability to approximate any target continuous func-
79 tion and classify any disjoint regions. Furthermore, their operation is fast and,
80 thus, they are appropriate for the cases where fast operation is important. How-
81 ever, most of the popular learning algorithms for SLFNs training are slow, due
82 to their iterative nature, and their parameter values should be carefully chosen.
83 This renders them inappropriate for dynamic classification schemes. Extreme

84 Learning Machine (ELM) [18] is a recently proposed algorithm for fast SLFNs
85 training requiring much less human effort. By using a sufficiently large number
86 of hidden neurons, the ELM classification scheme can be thought as a non linear
87 mapping of the training data in a high dimensional feature space, noted as ELM
88 space, followed by a linear classification procedure. Thus, non linear classifica-
89 tion functions can be approximated. Furthermore, the ELM training procedure
90 is independent of the training set size. These properties of ELM render it as a
91 good choice for dynamic classification schemes.

92 In this paper we propose a novel dynamic classification method inspired from
93 the above described dynamic subspace learning schemes and the effectiveness of
94 the ELM training procedure. The proposed classification procedure can be seen
95 as an adaptive multiple layer ANN, in which the number of layers, as well as the
96 number of each layer neurons, are dynamically determined by the test action in-
97 stance at hand, as illustrated in Figure 1. The proposed scheme is evaluated in ac-
98 tion recognition by using the dyname based action representation [19]. However,
99 it can be easily modified in order to be employed for different action representa-
100 tions. It is efficient in the sense that it dynamically determines the optimal labeled
101 data for training and classification. Furthermore, by exploiting the fast training
102 procedure of the ELM, the classification procedure is fast and efficient.

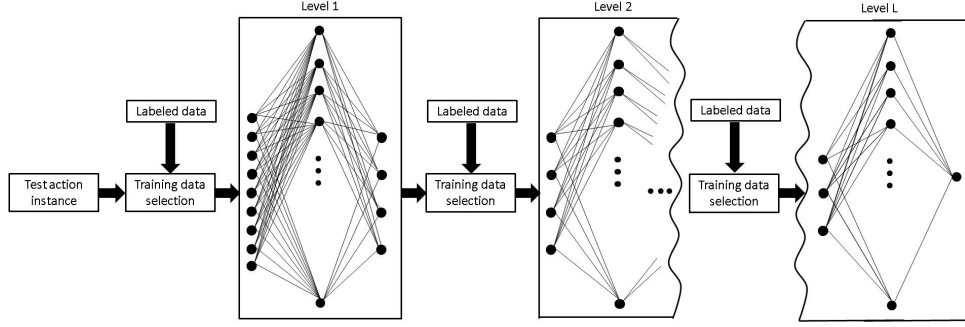


Figure 1: Adaptive multiple layer network topology.

103 The remainder of this paper is structured as follows. In Sections 2 we pro-
 104 vide an overview of the adopted action representation. In Sections 3 and 4 we
 105 present the two calculation steps that will be used in Section 5 to describe the pro-
 106 posed dynamic classification method. Section 6 presents experiments conducted
 107 for assessing its performance. Finally, conclusions are drawn in Section 7.

108 2. Dyneme based action representation

109 In this section, we present an overview of the dyneme based action represen-
 110 tation [19]. Let \mathcal{A} be an action class set consisting of C action classes, such as
 111 walk, run, jump, drink, eat, etc. Let \mathcal{U} denote an action recognition database
 112 containing N labeled action instances depicted in N videos, which will be called
 113 action videos hereafter. Video segmentation techniques, such as background sub-
 114 traction [20] or color based image segmentation [21], are applied to the action

115 video frames in order to produce binary action videos depicting the human body
 116 poses. The video frames forming the binary action videos are centered to the hu-
 117 man body regions of interest (ROIs), cropped to the ROIs region and resized to
 118 produce binary posture images of fixed ($H \times W$ pixels) size. In the experiments
 119 presented in this paper, we chose the size of the binary posture images to be equal
 120 to 32×32 pixels, which has been found experimentally to be a good compromise
 121 between computational cost and action recognition accuracy. The above described
 procedure is illustrated in Figure 2.

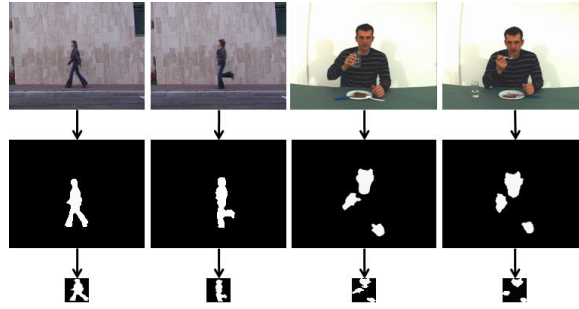


Figure 2: *Binary posture images production. From left to right 'walk', 'run', 'drink' and 'eat'.*

122

123 These binary images are represented as matrices, which are vectorized column-
 124 wise in order to produce the so called posture vectors \mathbf{p}_{ij} , $i = 1, \dots, N$, $j =$
 125 $1, \dots, N_i$, where N_i denotes the number of binary images forming binary action
 126 video i . Posture vectors of all the N labeled binary action videos are clustered,
 127 without exploiting the available label information, in order to produce D action

128 independent representative posture vectors, the dynemes. This is done by apply-
 129 ing D -Means clustering [11] to the posture vectors, minimizing the intra-cluster
 130 scatter, i.e.:

$$\sum_{d=1}^D \sum_{i=1}^N \sum_{j=1}^{N_i} \alpha_{ijd} \|\mathbf{p}_{ij} - \mathbf{v}_d\|^2, \quad (1)$$

131 where $\alpha_{ijd} = 1$, if \mathbf{p}_{ij} is assigned to cluster d and $\alpha_{ijd} = 0$, otherwise. Dynemes
 132 \mathbf{v}_d , $d = 1, \dots, D$ are defined to be the cluster mean vectors, i.e.:

$$\mathbf{v}_d = \frac{1}{n_d} \sum_{i=1}^N \sum_{j=1}^{N_i} \alpha_{ijd} \mathbf{p}_{ij}. \quad (2)$$

133 After dynemes calculation, each posture vector \mathbf{p}_{ij} is mapped to the mem-
 134 bership vector $\mathbf{u}_{ij} \in \mathbb{R}^D$, which denotes the fuzzy similarity of \mathbf{p}_{ij} with all the
 135 dynemes \mathbf{v}_d , according to a fuzzification parameter $m > 1$:

$$u_{ijd} = \frac{(\|\mathbf{p}_{ij} - \mathbf{v}_d\|_2)^{-\frac{2}{m-1}}}{(\sum_{k=1}^D \|\mathbf{p}_{ij} - \mathbf{v}_k\|_2)^{-\frac{2}{m-1}}}, \quad d = 1, \dots, D. \quad (3)$$

136 The optimal value of the fuzzification parameter m is obtained by applying
 137 the cross-validation procedure. Following [19], a value of $m = 1.1$ has been used
 138 in all the experiments presented in this paper. Finally, action vectors $\mathbf{s}_i \in \mathbb{R}^D$
 139 are calculated as the mean normalized membership vectors of the corresponding
 140 action videos:

$$\mathbf{s}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} \mathbf{u}_{ij}. \quad (4)$$

141 Action vectors \mathbf{s}_i representing all the training action videos are normalized in
 142 order to have zero mean and unit variance. Action vectors representing test action
 143 videos are normalized accordingly.

144 3. Data grouping and similarity measure

145 Let $\mathcal{Z} = \{\mathbf{z}_i\}_{i=1}^{N_z}$ be a vector set consisting of N_z labeled vectors \mathbf{z}_i . In order to
 146 determine K vector groups in \mathcal{Z} , we apply a clustering technique without exploit-
 147 ing the available action class labels. Since K -Means is a fast clustering algorithm,
 148 we employ K -Means to this end. That is, \mathcal{Z} is clustered by minimizing:

$$\sum_{k=1}^K \sum_{i=1}^{N_z} \beta_{ik} \|\mathbf{z}_i - \boldsymbol{\mu}_k\|^2. \quad (5)$$

149 $\boldsymbol{\mu}_k$ is the mean vector of group k , having cardinality $l_k = \sum_{i=1}^{l_k} \beta_{ik}$, i.e., $\boldsymbol{\mu}_k =$
 150 $\frac{1}{l_k} \sum_{i=1}^{l_k} \beta_{ik} \mathbf{z}_i$, and is used to represent the group. The number of groups K is
 151 either assumed to be known (fixed), or can be automatically determined. In the
 152 second case, several criteria can be used for optimal group number determination,
 153 such as the one described in [22].

154 In order to find the M most similar to a test vector \mathbf{z}_{test} vector groups, we
 155 calculate the Euclidean distances between \mathbf{z}_{test} and $\boldsymbol{\mu}_k$:

$$d_k = \|\mathbf{z}_{test} - \boldsymbol{\mu}_k\|^2. \quad (6)$$

After calculating d_k , $k = 1, \dots, K$, the M most similar vector groups to the test vector \mathbf{z}_{test} are those providing the M smallest distance values. M can either be assumed to be known (fixed), or can be automatically determined by following the procedure described in [15].

4. Extreme Learning Machine

Extreme Learning Machine (ELM) [18] is a fast algorithm for SLFNs training. In this section, we will provide an overview of ELM algorithm and discuss implementation issues appearing in our application setting. Let $\mathcal{X} = \{\mathbf{x}_i\}_1^{N_x}$ be a set of vectors, accompanied with the corresponding action class label set $\mathcal{C} = \{c_i\}_{i=1}^{N_x}$ $c_i \in \mathcal{A}$. The network's target vectors corresponding to each vector \mathbf{x}_i , $\mathbf{t}_i = [t_{i1}, \dots, t_{iC}]^T$, are set to $t_{ik} = 1$ for vectors belonging to action class k , i.e., when $c_i = k$, and $t_{ik} = -1$ otherwise.

In ELM, the network's input weights \mathbf{W}_{in} are randomly chosen, while the output weights \mathbf{W}_{out} are analytically calculated. Let us assume that the network's hidden layer consists of Q neurons and that $\mathbf{b} \in \mathbb{R}^Q$ is a vector containing the hidden layer neurons bias values, which are randomly chosen as well. Many activation functions $G()$ can be used for the hidden layer neurons' output calculation, such as sigmoid, sine, Gaussian and hard-limiting function. In our experi-

ments we have used the sigmoid function. That is, in our case $G(\mathbf{w}_j, b_j, \mathbf{x}_i) =$
 $\frac{1}{1+\exp^{-(\mathbf{w}_j^T \mathbf{x}_i + b_j)}}$, where \mathbf{w}_j denotes the j -th column of \mathbf{W}_{in} . By storing the hidden
layer neurons outputs in a matrix \mathbf{G} , i.e.,:

$$\mathbf{G} = \begin{bmatrix} G(\mathbf{w}_1, b_1, \mathbf{x}_1) & \cdots & G(\mathbf{w}_1, b_1, \mathbf{x}_{N_x}) \\ \vdots & \ddots & \vdots \\ G(\mathbf{w}_Q, b_Q, \mathbf{x}_1) & \cdots & G(\mathbf{w}_Q, b_Q, \mathbf{x}_{N_x}) \end{bmatrix}, \quad (7)$$

the network's output vector corresponding to the training vector \mathbf{x}_i can be written as $\mathbf{o}_i = \mathbf{W}_{out}^T \mathbf{g}_i$, where \mathbf{g}_i denotes the i -th column of \mathbf{G} . The network's outputs corresponding to the entire vector set \mathcal{X} can be written in a matrix form as $\mathbf{O} = \mathbf{W}_{out}^T \mathbf{G}$. Finally, by assuming that the network's predicted outputs \mathbf{O} are equal to the network's desired outputs \mathbf{T} , \mathbf{W}_{out} can be analytically calculated by $\mathbf{W}_{out} = \mathbf{G}^\dagger \mathbf{T}^T$, where $\mathbf{G}^\dagger = (\mathbf{G}\mathbf{G}^T)^{-1} \mathbf{G}$. However, the assumption of zero training error may decrease the generalization performance of the ELM network in the cases where the training set contains outliers. In order to increase the generalization performance of the ELM network, Huang et. al. [23] have recently proposed an optimization based regularized ELM algorithm formulated as follows:

$$\begin{aligned} \textbf{Minimize: } L_P &= \frac{1}{2} \|\mathbf{W}_{out}\|^2 + \Lambda \frac{1}{2} \sum_{i=1}^{N_x} \|\boldsymbol{\xi}_i\|^2 \\ \textbf{Subject to: } \mathbf{g}_i^T \mathbf{W}_{out} &= \mathbf{o}_i^T - \boldsymbol{\xi}_i^T, \quad i = 1, \dots, N_x, \end{aligned}$$

177 where $\xi_i \in \mathbb{R}^C$ is a training error vector corresponding to training sample \mathbf{x}_i and
 178 Λ is a parameter denoting the importance of the training error in the optimization
 179 problem. By adopting the above described optimization scheme, \mathbf{W}_{out} can be
 180 calculated by:

$$\mathbf{W}_{out} = \left(\frac{1}{\Lambda} \mathbf{I} + \mathbf{G}\mathbf{G}^T \right)^{-1} \mathbf{G}\mathbf{T}^T. \quad (8)$$

181 After \mathbf{W}_{out} calculation, a test vector \mathbf{x}_{test} can be introduced to the ELM net-
 182 work and be classified to the class corresponding to the highest network's output,
 183 i.e.:

$$c_{test} = \arg \max_j o_{test,j}, j = 1, \dots, C. \quad (9)$$

184 As can be seen the ELM training procedure is fast, since it involves matrix
 185 multiplication and matrix inversion operations. Such operations can be efficiently
 186 calculated by existing optimized software [24, 25]. Furthermore, the network
 187 topology and the input weights \mathbf{W}_{in} can be determined only once, since they do
 188 not involve any training procedure.

189 **5. Dynamic classification scheme**

190 In this section we present the proposed dynamic classification method. Let
 191 \mathcal{U} be an action recognition database, containing N action videos accompanied by
 192 the corresponding action class labels $c_i, i = 1, \dots, N$ belonging to C action classes

193 forming an action class set \mathcal{A} . These action videos are preprocessed, following
 194 the procedure described in Section 2, in order to produce N action vectors $\mathbf{s}_i \in$
 195 \mathbb{R}^D , $i = 1, \dots, N$.

196 Most classification schemes would employ all the available labeled action vec-
 197 tors \mathbf{s}_i , $i = 1, \dots, N$ and the corresponding action class labels c_i in order to cal-
 198 culate a static classification model, that would be used in order to classify any
 199 unknown (test) action vector. In our case, the set of action vectors $\mathcal{S} = \{\mathbf{s}_i\}_{i=1}^N$
 200 is clustered, by performing the procedure described in Section 3, in order to de-
 201 termine K action vector groups, represented by the corresponding mean group
 202 vectors $\boldsymbol{\mu}_k \in \mathbb{R}^D$, $k = 1, \dots, K$.

203 Let a test action video be represented by an action vector $\mathbf{s}_{test} \in \mathbb{R}^D$. \mathbf{s}_{test}
 204 is compared with all the K mean group vectors $\boldsymbol{\mu}_k$ in order to determine the M
 205 closest to \mathbf{s}_{test} groups. The action vectors belonging to these M groups form the
 206 algorithm's first level training set $\mathcal{S}_1 = \{\mathbf{s}_{i,1}\}_{i=1}^{N_1}$. Here we have introduced a
 207 second index denoting the levels of the proposed dynamic classification scheme.
 208 Action class labels corresponding to the action vectors forming \mathcal{S}_1 are employed
 209 in order to form the first level action class label set $\mathcal{C}_1 = \{c_{i,1}\}_{i=1}^{N_1}$, $c_{i,1} \in \mathcal{A}_1$.
 210 Obviously, $\mathcal{A}_1 \subseteq \mathcal{A}$, since only the labeled action vectors belonging to the action
 211 classes that are most similar to the actual \mathbf{s}_{test} action class are included in \mathcal{S}_1 . Now,

212 we can formulate an alternative classification problem. Instead of employing the
 213 entire action vector set \mathcal{S} and train a universal classifier, we can use the action
 214 vector set \mathcal{S}_1 in order to train a \mathbf{s}_{test} - specific classifier. That is, we train an SLFN
 215 by using \mathcal{S}_1 and \mathcal{C}_1 following the procedure described in Section 4. Subsequently,
 216 we introduce \mathbf{s}_{test} to the trained network and we obtain its response \mathbf{o}_{test} . In this
 217 stage, we can classify \mathbf{s}_{test} to the action class that provides the maximal network
 218 output, i.e.:

$$c_{test} = \underset{j}{argmax} o_{test,j}. \quad (10)$$

219 However, we choose to perform the dynamic classification procedure in mul-
 220 tiple levels L . For this reason, we introduce \mathcal{S}_1 to the trained network and we
 221 obtain its responses $\mathcal{O}_1 = \{\mathbf{o}_{i,1}\}_{i=1}^{N_1}$. By using \mathcal{O}_1 , we can now reformulate the
 222 classification problem. In the general case, after obtaining the l -th level network
 223 outputs, $\mathcal{O}_l = \{\mathbf{o}_{i,l}\}_{i=1}^{N_l}$ and $\mathbf{o}_{test,l}$, the feature vectors forming \mathcal{O}_l are grouped by
 224 following the procedure described in Section 3. $\mathbf{o}_{test,l}$ is, subsequently, compared
 225 with the corresponding mean group vectors $\boldsymbol{\mu}_{k,l}$ and the closest to $\mathbf{o}_{test,l}$ groups
 226 are used to form the $(l+1)$ -th level training set \mathcal{S}_{l+1} . The $(l+1)$ -th level network
 227 is, subsequently, trained by using \mathcal{S}_{l+1} and the corresponding action class label
 228 set $\mathcal{C}_{l+1} = \{c_{l+1,i}\}_{i=1}^{N_{l+1}}$, $c_{l+1,i} \in \mathcal{A}_{l+1}$. Obviously the number of action classes
 229 forming the classification problem of every level of the proposed dynamic clas-

230 sification scheme is a subset of the number of classes forming the previous level
 231 classification problem, i.e., $\mathcal{A}_{l+1} \subseteq \mathcal{A}_l$.

232 The above described iterative procedure is performed multiple times, until the
 233 vectors forming the network training set belong to one action class only. That is,
 234 the maximal number of classification levels L depends on the test action vector
 235 \mathbf{s}_{test} . In the cases where the action class that \mathbf{s}_{test} belongs to, is well distinguished
 236 from all the other action classes forming \mathcal{A} , only one classification level will be
 237 performed. In the cases of overlapping action classes, multiple classification levels
 238 will be performed in order to obtain the final classification result. Since at each
 239 level of the dynamic classification procedure the network training set is a subset of
 240 the previous level network training set, i.e., $\mathcal{S}_l \subset \mathcal{S}_{l-1}$, and the number of available
 241 labeled action vectors is finite, the proposed iterative procedure will converge in
 242 a finite number of iterations. In the, extreme, case of highly overlapping action
 243 classes, the iterative procedure will end when the network training set consists of
 244 only one labeled vector.

245 Consider the example illustrated in Figure 3. In this Figure, we illustrate the
 246 2-dimensional feature space resulted by applying Principal Component Analysis
 247 (PCA) [11] on the dyneme based action video representation, in the Weizemann
 248 action recognition database [26], which will be used in the first set of the experi-

ments presented in the following section.

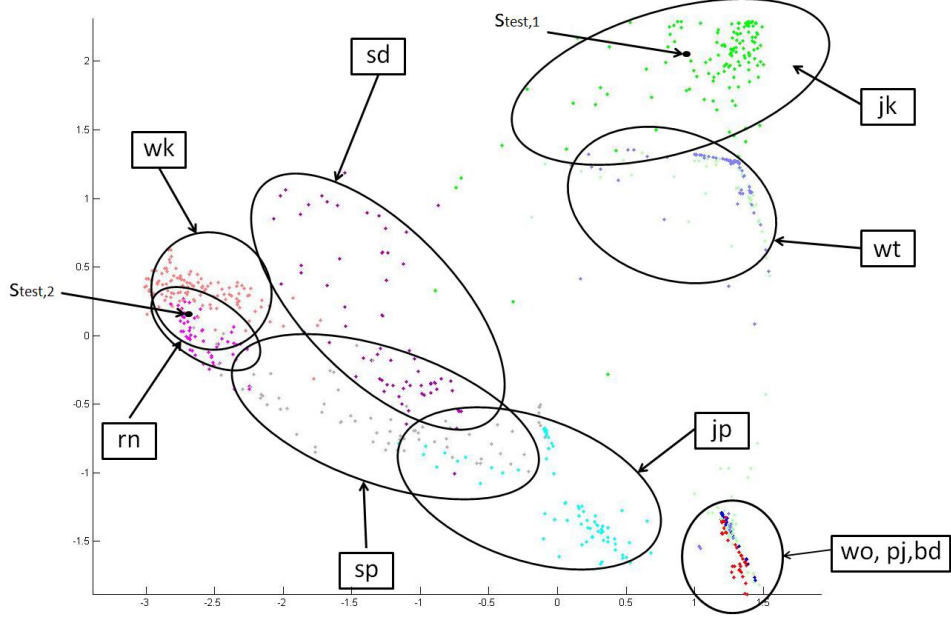


Figure 3: 2D space resulted by applying PCA on the dyneme based action representation in the Weizmann action recognition database.

249

250 As can be seen in Figure 3, some action classes, such as 'jumping jack' (jk),
 251 may be well distinguished from all the other action classes. However, action
 252 classes, usually, are confused with each other. Similar action classes, such as
 253 'walk' (wk) and 'run' (rn), or 'skip' (sp) and 'jump in place' (jp) contain a high
 254 number of common human body poses and, thus, variations in action execution
 255 style and human body size may result to similar action representations. Assume
 256 that a test action video, represented by the corresponding action vector $s_{test,1}$, be-

257 longs to the action class 'jumping jack' (jk). In this case, it is expected that $s_{test,1}$
258 will be directly classified to the correct action class, since 'jumping jack' is well
259 distinguished from all the other action classes. However, in the case of a test ac-
260 tion video belonging to the action class 'run' (rn), represented by the action vector
261 $s_{test,2}$, its classification procedure is not obvious, since action class 'run' is con-
262 fused with action class 'walk'. Thus, in this case, the classification procedure will
263 probably involve multiple classification levels.

264 In Figure 3, it can also be seen that action classes are not linearly separable.
265 For example, consider the case of action classes 'walk' and 'run', as highlighted in
266 Figure 4. Clearly, these two action classes share the same feature space and are not
267 linearly separable. Thus, the use of linear models for action class discrimination
268 is not an appropriate choice. This can be seen in Figure 5, where we illustrate the
269 separating hyperplanes (lines) resulted by applying LDA (Figure 5a) and ELM
270 (Figure 5b) based action vectors classification, respectively. It can be seen that
271 by applying the ELM based action vector classification, action classes are better
272 discriminated. This is reasonable, since, as it was previously discussed, the use of
273 ELM can better capture the non linear nature of the action classes.

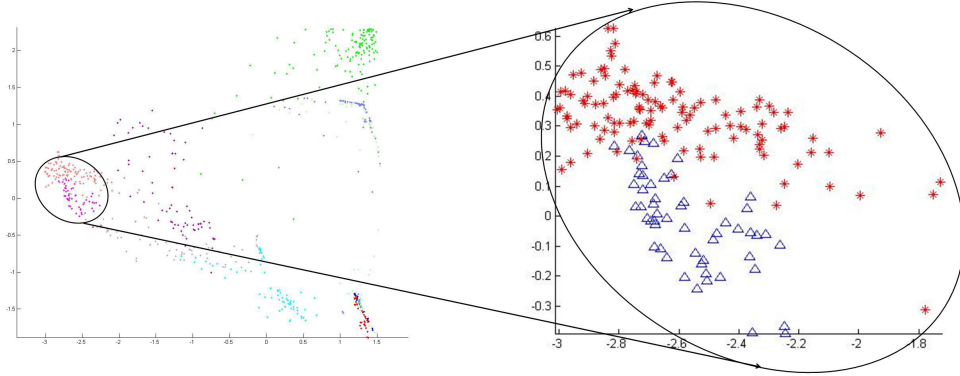


Figure 4: Action vectors belonging to action classes 'walk' and 'run' in the 2D PCA space.

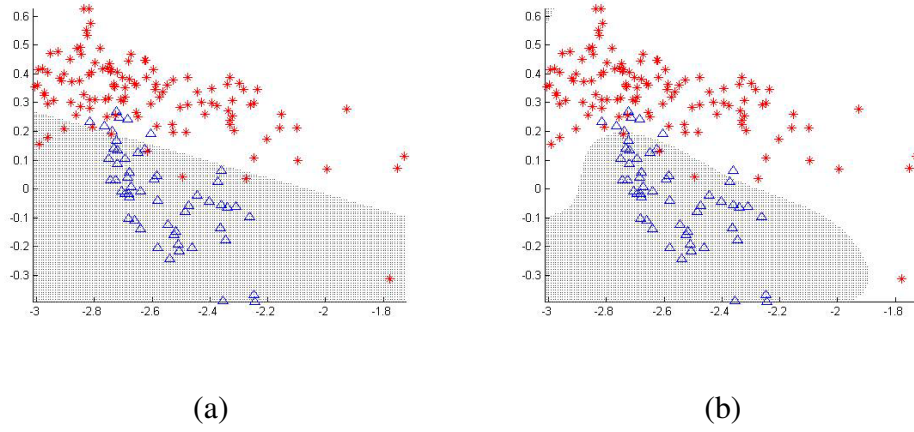


Figure 5: Separating planes (lines) resulted by: a) applying LDA followed by nearest class centroid and b) training a single hidden layer neural network using the ELM algorithm.

274 6. Experimental results

275 In this Section we present experiments conducted in order to evaluate the pro-
 276 posed dynamic classification method. We conducted experiments on the Weize-
 277 mann [26] and the i3DPost [27] action recognition databases containing daily ac-

278 tions, as well as on a new action recognition database aiming at recognition of
 279 actions appearing in meal intakes [28]. We provide a comprehensive descrip-
 280 tion of these databases in Subsections 6.1, 6.3 and 6.5, respectively. In each
 281 level of the proposed dynamic classification scheme, we grouped the labeled
 282 vectors in $K = [10, 20, 50]$ groups. The optimal number of closest to the test
 283 vector groups has been experimentally determined by using different values of
 284 $M = \frac{K}{k}$, $k = 1, \dots, 10$. Regarding the optimal number of dynemes D , the num-
 285 ber of network hidden layer neurons Q and the parameter value Λ of the ELM
 286 algorithm, they have been determined by performing the leave-one-out cross-
 287 validation procedure. Specifically, we have performed the cross-validation pro-
 288 cedure using values of D equal to $10k$, $k = 1, \dots, 20$, $Q = [100, 200, 500, 1000]$
 289 and values of Λ equal to 10^λ , $\lambda = -5, \dots, 5$. In order to assess the ability of the
 290 proposed classification scheme to generalize on data that it was not trained on, we
 291 performed the leave-one-person-out cross-validation procedure (LOPOCV). That
 292 is, we used the action videos depicting all but one person in the database as la-
 293 beled data and the action videos depicting the remaining one as test data, in order
 294 to perform one iteration (fold) of the cross validation procedure. Multiple folds,
 295 equal to the number of persons appearing in the database, have been performed in
 296 order to complete an experiment.

297 6.1. *Weizemann database*

298 The Weizemann action recognition database [26] contains 90 low-resolution,
299 144 × 180 pixel, image sequences depicting nine persons (five males and four
300 females) performing ten daily actions each. The actions appearing in the database
301 are: 'walk' (wk), 'run' (rn), 'jump in place on two legs' (pj), 'jump forward on two
302 legs' (jp), 'jumping-jack' (jk), 'gallop sideways' (sd), 'skip' (sp), 'wave one hand'
303 (wo), 'wave two hands' (wt) and 'bend' (bd). Binary image sequences denoting
304 the human body regions are included in the database. Example video frames and
305 binary skin-colored regions are illustrated in Figure 2.

306 Since most of these image sequences depict multiple action instances, e.g.
307 multiple walking steps, we automatically produced binary action videos by using
308 the binary image sequences and a sliding window consisting of 16 video frames,
309 moving in steps of 4 video frames, resulting to the creation of 952 action videos.
310 Figure 6 illustrates the sliding window technique for automatic binary action video
311 creation. The resulted binary action videos have been preprocessed following the
312 procedure described in Section 2

313 6.2. *Experiments on the Weizemann database*

314 In our first set of experiments we have conducted the LOPOCV procedure
315 on the Weizemann action recognition database using the resulted binary action

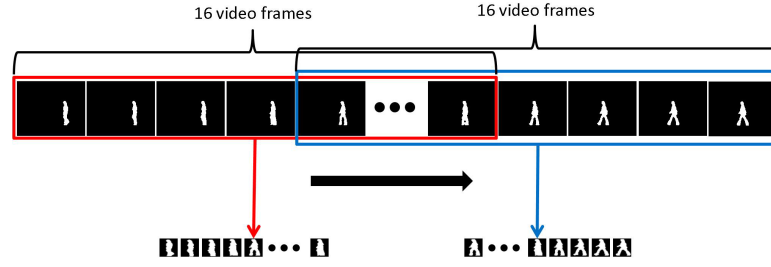


Figure 6: *Sliding window technique for automatic action video creation.*

316 videos. In Figure 7 we illustrate the action classification rates obtained by us-
 317 ing different values of M . As can be seen, by using smaller values of M the
 318 classification rate increases. This is reasonable since for smaller values of M the
 319 classification procedure involves only the labeled data that are more similar to the
 test ones.

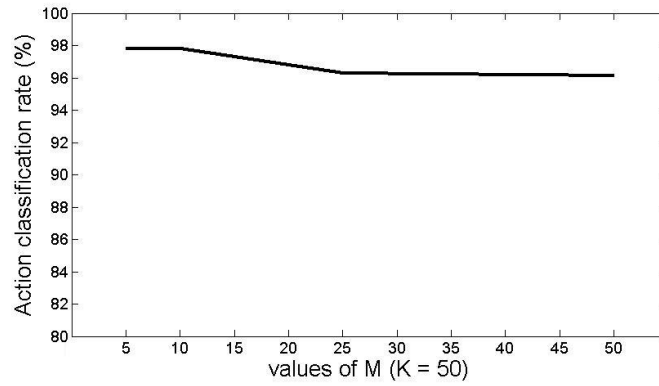


Figure 7: *Action classification rates as a function of M .*

320

321 The confusion matrix corresponding to the optimal parameters is illustrated in
 322 Figure 8a. As can be seen, high classification rates have been obtained for all the

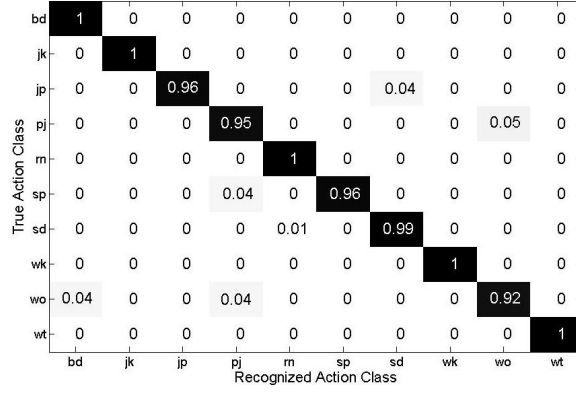
323 action classes. The class which was found to be the most difficult for classification
324 is action class 'wave one hand', which is confused with action classes 'bend' and
325 'jump in place on two legs'. However, even for this case a high classification rate,
326 equal to 92%, has been obtained.

327 In order to directly compare the performance of the proposed classification
328 method with other ones, we have conducted experiments by performing the LOPOCV
329 procedure on the Weizemann action recognition database using both static and
330 dynamic classification strategies. That is, we performed the LOPOCV procedure
331 by employing the static classification strategy and performing LDA based action
332 vector projection followed by nearest class centroid classification, resulting to an
333 action classification rate equal to 95.92%. By following the static classification
334 strategy and performing ELM based action vector classification, an action classi-
335 fication rate equal to 96.15% has been obtained.

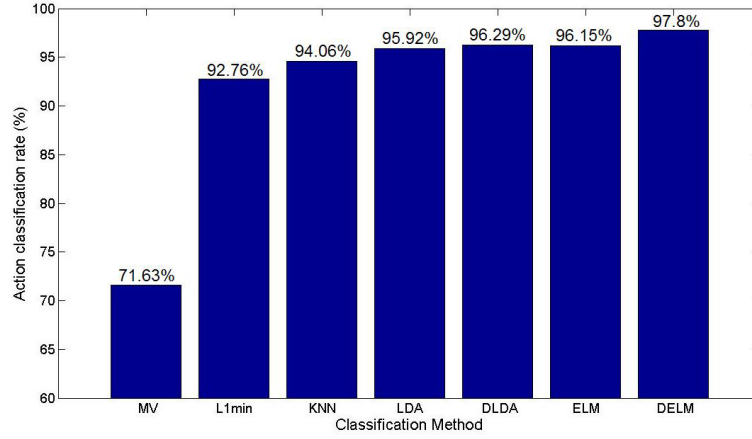
336 Subsequently, we have conducted experiments employing the dynamic classi-
337 fication strategy. KNN action vectors classification, using $K = 3$ nearest neigh-
338 bors, resulted to an action classification rate equal to 94.06%. Action classification
339 based on L1-minimization followed by smallest residual error action vector clas-
340 sification, as proposed by Wright et. al. [13], resulted to an action classification
341 rate equal to 92.76%. By following the dynamic classification method proposed

342 by Kyperountas et. al. [15], which employs LDA based data projection, an action
343 classification rate equal to 96.29% has been obtained. Finally, by performing one
344 level of the proposed dynamic classification method and classifying each test ac-
345 tion vector by applying majority voting on the action class labels of the labeled
346 action vectors forming the M , out of K , closest to the test action vector groups, an
347 action classification rate equal to 71.63% has been obtained. This procedure can
348 be used as a reference for the performance of the proposed dynamic classification
349 scheme, since, intuitively, the determination of labeled action vectors similar to
350 the test one should lead to correct classification results. The action classification
351 rates obtained in all these experiments are summarized in Figure 8b.

352 As can be seen in Figure 8, the adoption of a dynamic classification strategy
353 leads to an increase of the action classification rates. In both the LDA and ELM
354 cases, the dynamic classification approach provides higher classification rates.
355 Furthermore, it can be seen that the proposed classification scheme is efficient,
356 since a simple majority voting on the action labels of the labeled action vectors
357 that form the M most similar to the test action vectors groups, results to a, rel-
358 atively, high action classification rate. Finally, it can be seen that the proposed
359 dynamic classification scheme outperforms all the other competing methods ap-
360 pearing in Figure 8, since it combines both efficient search of the most appropriate



(a)



(b)

Figure 8: a) Confusion matrix on the Weizemann database obtained by applying the proposed dynamic action classification method and b) Comparison results on the Weizemann action recognition database.

361 training set and approximation of non-linear discrimination functions.

362 In order to compare the performance of the proposed action classification

363 method with that of other methods proposed in the literature, we have followed

the procedure proposed in [29]. That is, the image sequences of the Weizemann database have been classified to action classes by performing majority voting on the action classification results provided by the algorithm for the corresponding action videos, resulting to an action classification rate equal to 98.9%. All but one action sequences have been correctly classified. The only sequence that was misclassified belongs to action class 'skip' and classified to action class 'jump forward on two legs'. Since some of the methods proposed in the literature providing state of the art performance are evaluated by using an earlier version of the database containing nine action classes, i.e., not containing action class 'skip', we have also tested the proposed dynamic action classification method by using this earlier version. Comparison results with other action recognition methods are illustrated in Table 1.

6.3. AIIA-MOBISERV database

Despite the fact that most applications, including action recognition functionality, consider daily action types, such as walk, run, etc., there are applications requiring different type of actions. For example, monitoring the status of the elderly people in the early stages of dementia, while still living independently, to prevent dehydration is an important task. In the framework of the EU R&D project MOBISERV, we created an eating and drinking action recognition database, which

Table 1: Comparison results on the Weizmann action recognition database.

Method	9 actions	10 actions
Yaffet & Wolf [30]	100%	–
Wang & Mori [31]	100%	–
Guha & Ward [32]	–	98.9%
Gorelic et al. [29]	–	97.8%
Riemenchneider et al. [33]	–	96.7%
Gkalelis et.al. [34]	–	96%
Ali & Shah [35]	–	95.7%
Junejo et al. [36]	95.3%	–
Thureau & Hlavac [37]	–	94.4%
Zhang et al. [38]	–	92.8%
Niebles et al. [39]	–	90%
Proposed method	100%	98.9%

383 is publicly available in [28]. Twelve persons (six females and six males) were
 384 captured by a camera placed at a distance of 2 meters in front of them, during a
 385 meal. Four meals have been recorded, each for a different day for all the twelve
 386 persons. The actions appearing in the database are: 'eat', 'drink' and 'apraxia'.
 387 Action class 'eat' contains the cases where the person eats using a spoon, a cut-

388 lery, a fork, or takes a bite using one or two hands. Action class 'drink' contains
389 the cases where the person drinks using a cup, a glass, or a straw. Finally, action
390 class 'apraxia' contains the cases where the person is slicing his/her food or he/she
391 is chewing it and the cases where the person rests.

392 We have manually temporally segmented the videos depicting all the persons
393 during two meals. This procedure resulted to the creation of 1288 action videos.
394 A color based image segmentation technique has been applied to the video frames
395 of these action videos in order to produce binary images depicting the skin regions
396 of the depicted person's body. Specifically, each video frame has been converted
397 to the HSV color space and the image pixels having HS values in pre-specified
398 thresholds, corresponding to skin-like color values, have been determined to be
399 foreground pixels, while the rest pixels have been assumed to belong to the back-
400 ground. Morphological operations (closing) have been, subsequently, performed
401 in order to obtain the final binary action video frames. This resulted to the cre-
402 ation of binary action videos denoting the person's head and hands, which have
403 been preprocessed following the procedure described in Section 2. Example video
404 frames and binary skin-colored regions are illustrated in Figure 2.

405 6.4. Experiments on the AIIA-MOBISERV database

406 In our second set of experiments we have performed the LOPOCV procedure
407 on the binary action videos of the AIIA-MOBISERV database. An action classi-
408 fication rate equal to 93.4% has been obtained by applying the proposed dynamic
409 action classification method. The confusion matrix of this experiment is illus-
410 trated in Figure 9a. Comparison results with other dynamic, as well as static,
411 action classification schemes are illustrated in Figure 9b. As can be seen, by ap-
412 plying the majority voting classification scheme, an action classification rate equal
413 to 65.11% has been obtained. Action classification based on L1-minimization fol-
414 lowed by smallest residual error action vector classification, resulted to an action
415 classification rate equal to 90.3%. KNN ($K = 3$) action vector classification re-
416 sulted to an action classification rate equal to 89.91%. Static LDA and ELM based
417 action classification schemes, provided action classification rates equal to 89.94%
418 and 89.73%, respectively. Finally, LDA-based dynamic action vector classifica-
419 tion resulted to an action classification rate equal to 92.53%. As can be seen, the
420 dynamic action classification approach outperforms the static one in both the LDA
421 and ELM based classification schemes. Furthermore, it can be seen that the ELM
422 based dynamic action classification scheme outperforms all the methods presented
423 in this Figure.

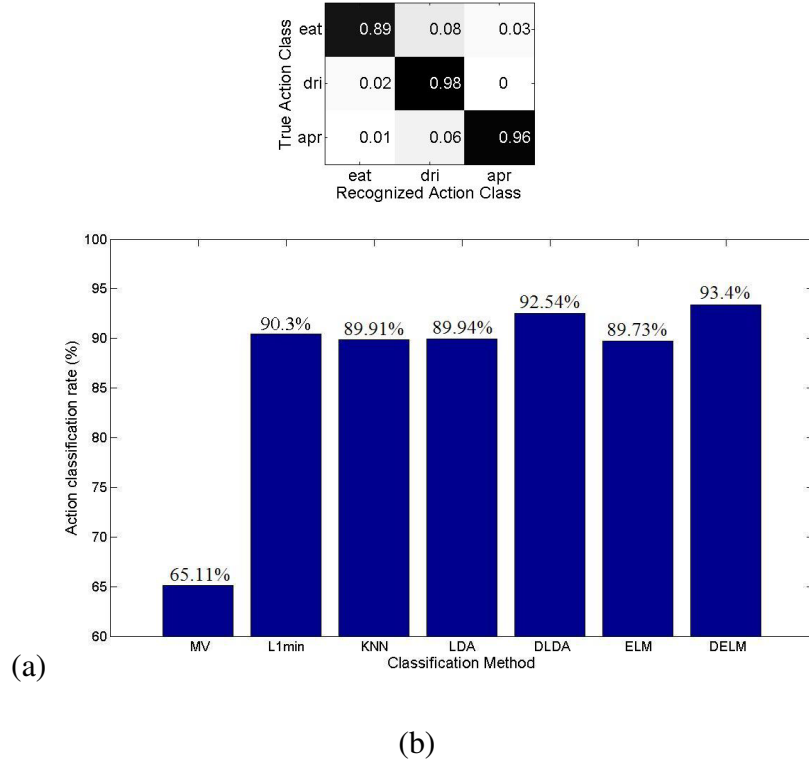


Figure 9: a) Confusion matrix on the AIIA-MOBISERV database obtained by applying the proposed dynamic action classification method and b) Comparison results on the AIIA-MOBISERV action recognition database.

6.5. i3DPost database

The i3DPost multi-view database [27] contains 512 high resolution (1080×1920 pixels) image sequences depicting eight persons (six males and two females) performing eight actions. The database camera setup consists of eight cameras, providing a 360° coverage of the scene. The actions appearing in the database are: 'walk' (wk), 'run' (rn), 'jump in place' (jp), 'jump forward' (jf), 'bend' (bd), 'fall

430 down' (fl), 'sit on a chair' (st) and 'wave one hand' (wo). Since most of the image
431 sequences depict multiple action instances, e.g. multiple walking steps, we have
432 manually temporally segmented them in order to produce videos depicting one
433 action instance each. A color based image segmentation technique, discarding the
434 blue color in the HSV color space, has been applied to the video frames of these
435 action videos in order to produce binary action videos denoting the human body.
436 Morphological operations (closing) have been, subsequently, performed in order
437 to obtain the final binary action video frames.

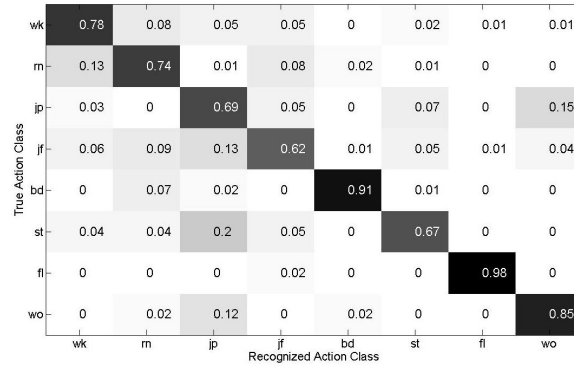
438 6.6. *Experiments on the i3DPost database*

439 In our third set of experiments we have performed the LOPOCV procedure on
440 the binary action videos of the i3DPost database. In each fold of the LOPOCV
441 procedure, we have used the action videos depicting seven of the persons perform-
442 ing an action instance from all the available cameras as labeled data. Each action
443 video depicting the test person has been classified to one of the eight action classes
444 independently, in order to form a single-view view-invariant action classification
445 problem. We should note that, we expected the above described procedure to re-
446 sult to a difficult classification problem due to the well known view angle effect
447 [40]. By applying the proposed dynamic action classification method, an action
448 classification rate equal to 77.97% has been obtained. The confusion matrix of this

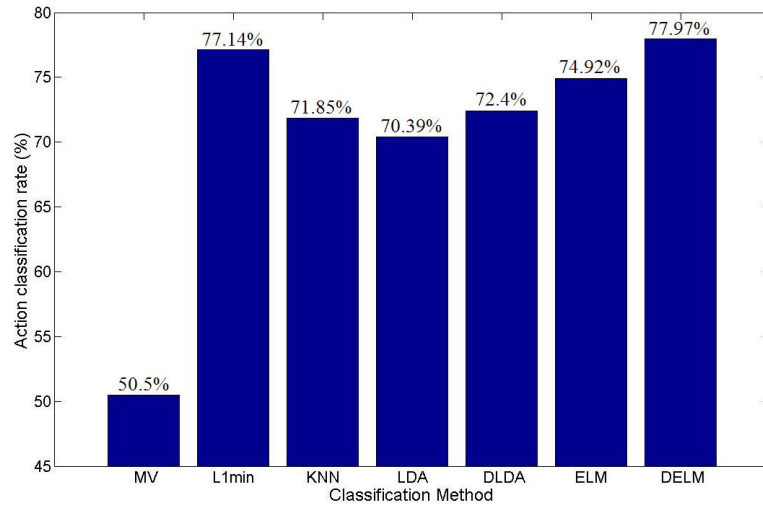
experiment is illustrated in Figure 10a. Comparison results with other dynamic, as well as static, action classification schemes are illustrated in Figure 9b. As can be seen, by applying the majority voting classification scheme, an action classification rate equal to 50.5% has been obtained. Action classification based on L1-minimization followed by smallest residual error action vector classification, resulted to an action classification rate equal to 77.14%. KNN ($K = 3$) action vector classification resulted to an action classification rate equal to 71.85%. Static LDA and ELM based action classification schemes, provided action classification rates equal to 70.39%. Finally, LDA-based dynamic action vector classification resulted to an action classification rate equal to 72.4%. As can be seen, the dynamic action classification approach outperforms the static one in both the LDA and ELM based classification schemes. Furthermore, it can be seen that the ELM based dynamic action classification scheme outperforms all the methods presented in this Figure.

7. Conclusions

In this paper, we proposed a novel dynamic action classification method based on an iterative procedure determining test action instance specific classification problems in multiple levels. Action instances are represented by vectors denoting



(a)



(b)

Figure 10: a) Video frames depicting one person of the i3DPost database walking from different viewing angles and b) Comparison results on the i3DPost action recognition database.

467 the fuzzy similarity of the corresponding human body poses with representative
468 human body poses, the dynemes. At each classification level, the most similar
469 to the test action instance labeled vectors are employed in order to train a single

470 hidden layer feedforward network using the ELM algorithm. A new feature space
471 is, subsequently, obtained by the trained network's outputs. By exploiting the
472 properties of the adopted network topology and the fast training procedure of the
473 ELM algorithm, the proposed classification method is fast and efficient. Experi-
474 ments on publicly available action recognition databases indicate the superiority
475 of the dynamic classification strategy, compared to the static one, as well as the
476 effectiveness of the proposed dynamic classification method.

477 **Acknowledgment**

478 The research leading to these results has received funding from the Collabo-
479 rative European Project MOBISERV FP7-248434 (<http://www.mobiserv.eu>), An
480 Integrated Intelligent Home Environment for the Provision of Health, Nutrition
481 and Mobility Services to the Elderly.

482 **References**

- 483 [1] L. Weilun, H. Jungong, P. With, Flexible Human Behavior Analysis Framework
484 for Video Surveillance Applications, International Journal of Digital Multimedia
485 Broadcasting 2010 (2010) 9, ISSN 1687-7578.
- 486 [2] P. Barr, J. Noble, R. Biddle, Video game values: Human-computer interaction and
487 games, Interacting with Computers 19 (2) (2007) 180–195.

- 488 [3] T. Hollerer, S. Feiner, D. Hallaway, B. Bell, M. Lanzagorta, D. Brown, S. Julier,
489 Y. Baillet, L. Rosenblum, User interface management techniques for collaborative
490 mobile augmented reality, *Computers and Graphics* 25 (5) (2001) 799–810.
- 491 [4] J. Assfalg, M. Bertini, C. Colombo, A. Del Bimbo, W. Nunziati, Semantic anno-
492 tation of soccer videos: automatic highlights identification, *Computer Vision and*
493 *Image Understanding* 92 (2-3) (2003) 285–305.
- 494 [5] S. Ali, M. Shah, Human Action Recognition in Videos Using Kinematic Features
495 and Multiple Instance Learning, *IEEE Transactions on Pattern Analysis and Ma-*
496 *chine Intelligence* 32 (2) (2010) 288–303.
- 497 [6] Y. Wang, G. Mori, Hidden Part Models for Human Action Recognition: Proba-
498 bilistic versus Max Margin, *IEEE Transactions on Pattern Analysis and Machine*
499 *Intelligence* 33 (7) (2011) 1310–1323.
- 500 [7] M. Giese, T. Poggio, Neural mechanisms for the recognition of biological move-
501 ments, *Nature Reviews Neuroscience* 4 (3) (2003) 179–192.
- 502 [8] A. Iosifidis, A. Tefas, N. Nikolaidis, I. Pitas, Multi-view human movement recogni-
503 tion based on Fuzzy distances and Linear Discriminant Analysis, *Elsevier Computer*
504 *Vision and Image Understanding* 116 (2012) 347–360.
- 505 [9] A. Iosifidis, A. Tefas, I. Pitas, View-Invariant Action Recognition Based on Artificial

- 506 Neural Networks, *IEEE Transactions on Neural Networks and Learning Systems*
507 23 (3) (2012) 412–425.
- 508 [10] S. Haykin, *Neural networks and learning machines*, Upper Saddle River, New Jersey
509 .
- 510 [11] R. Duda, P. Hart, D. Stork, *Pattern Classification*, 2nd ed, Wiley-Interscience, 2000.
- 511 [12] A. Iosifidis, A. Tefas, I. Pitas, Person specific activity recognition using fuzzy learn-
512 ing and discriminant analysis, *European Signal Processing Conference (EUSIPCO)*
513 (2011) 1974–1978.
- 514 [13] J. Wright, A. Yang, A. Ganesh, S. Sastry, Y. Ma, Robust face recognition via sparse
515 representation, *IEEE Transactions on Pattern Analysis and Machine Intelligence*
516 31 (2) (2009) 210–227.
- 517 [14] H. Tang, M. Lyu, I. King, Face recognition committee machines: dynamic vs. static
518 structures, *Proceedings of 12th International Conference on Image Analysis and*
519 *Processing* (2003) 121–126.
- 520 [15] M. Kyperountas, A. Tefas, I. Pitas, Dynamic training using multistage clustering for
521 face recognition, *Pattern Recognition* 41 (3) (2008) 894–905.
- 522 [16] J. Lu, K. Plataniotis, A. Venetsanopoulos, Face recognition using LDA-based algo-
523 rithms, *IEEE Transactions on Neural Networks* 14 (1) (2003) 195–200.

- 524 [17] J. Lu, K. Plataniotis, A.N.Venetsanopoulos, Regularization studies of linear discrim-
525 inant analysis in small sample size scenarios with application to face recognition,
526 Pattern Recognition Letters 26 (2) (2005) 181–191.
- 527 [18] G. Huang, Q. Zhu, C. Siew, Extreme learning machine: a new learning scheme of
528 feedforward neural networks, IEEE International Joint Conference on Neural Net-
529 works 2 (2004) 985–990.
- 530 [19] N. Gkalelis, A. Tefas, I. Pitas, Combining fuzzy vector quantization with linear dis-
531 criminant analysis for continuous human movement recognition, IEEE Transactions
532 on Circuits and Systems for Video Technology 18 (11) (2008) 1511–1521.
- 533 [20] M. Seki, T. Wada, H. Fujiwara, K. Sumi, Background subtraction based on the cooc-
534 currence of image variations, Proceedings of IEEE Conference on Computer Vision
535 and Pattern Recognition 2 (2003) 65–72.
- 536 [21] V. Vezhnevets, V. Sazonov, A. Andreeva, A survey on pixel-based skin color detec-
537 tion techniques, Proc. Graphicon 3.
- 538 [22] B. Frey, D. Dueck, Clustering by passing messages between data points, science
539 315 (5814) (2007) 972.
- 540 [23] G. Huang, H. Zhou, X. Ding, R. Zhang, Extreme Learning Machine for Regression

- 541 and Multiclass Classification, IEEE Transactions on Systems, Man, and Cybernet-
542 ics, Part B: Cybernetics 42 (2) (2012) 513–529.
- 543 [24] G. Bradski, A. Kaehler, Learning OpenCV: Computer vision with the OpenCV li-
544 brary, O'Reilly Media, 2008.
- 545 [25] <http://icl.cs.utk.edu/lapack-for-windows/> .
- 546 [26] M. Blank, L. Gorelick, E. Shechtman, M. Irani, R. Basri, Actions as space-time
547 shapes, Proceedings of IEEE Conference on Computer Vision 2 (2005) 1395–1402.
- 548 [27] N. Gkalelis, H. Kim, A. Hilton, N. Nikolaidis, I. Pitas, The i3DPost multi-view and
549 3D human action/interaction database, 6th Conference on Visual Media Production
550 (2009) 159–168.
- 551 [28] <http://poseidon.csd.auth.gr/MOBISERV-AIIA/index.html> .
- 552 [29] L. Gorelick, M. Blank, E. Shechtman, M. Irani, R. Basri, Actions as Space-Time
553 Shapes, Transactions on Pattern Analysis and Machine Intelligence 29 (12) (2007)
554 2247–2253.
- 555 [30] L. Yeffet, L. Wolf, Local trinary patterns for human action recognition, International
556 Conference on Computer Vision (2009) 492–497.
- 557 [31] Y. Wang, G. Mori, Human action recognition by semilattent topic models, IEEE

- 558 Transactions on Pattern Analysis and Machine Intelligence 31 (10) (2009) 1762–
559 1774.
- 560 [32] T. Guha, R. Ward, Learning Sparse Representations for Human Action Recognition,
561 IEEE Transactions on Pattern Analysis and Machine Intelligence .
- 562 [33] H. Riemenschneider, M. Donoser, H. Bischof, Bag of optical flow volumes for image
563 sequence recognition, British Machine Vision Conference (BMVC) .
- 564 [34] N. Gkalelis, A. Tefas, I. Pitas, Combining fuzzy vector quantization with linear dis-
565 criminant analysis for continuous human movement recognition, IEEE Transactions
566 on Circuits and Systems for Video Technology 18 (11) (2008) 1511–1521, ISSN
567 1051-8215.
- 568 [35] S. Ali., M. Shah, Human action recognition in videos using kinematic features and
569 multiple instance learning, IEEE Transactions on Pattern Analysis and Machine In-
570 telligence 32 (2) (2010) 288–303.
- 571 [36] I. Junejo, E. Dexter, I. Laptev, P. Pérez, View-independent action recognition from
572 temporal self-similarities, IEEE Transactions on Pattern Analysis and Machine In-
573 telligence 33 (1) (2011) 172–185.
- 574 [37] C. Thureau, V. Hlavác, Pose primitive based human action recognition in videos or

- 575 still images, Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE
576 Conference on (2008) 1–8.
- 577 [38] Z. Zhang, Y. Hu, S. Chan, L. Chia, Motion context: A new representation for human
578 action recognition, European Conference on Computer Vision ECCV (2008) 817–
579 829.
- 580 [39] J. Niebles, H. Wang, L. Fei-Fei, Unsupervised learning of human action cate-
581 gories using spatial-temporal words, International Journal of Computer Vision 79 (3)
582 (2008) 299–318.
- 583 [40] S. Yu, D. Tan, T. Tan, Modeling the effect of view angle variation on appearance-
584 based gait recognition, Proceedings Asian Conf. Computer Vision 1 (2006) 807–
585 816.